

به نام خدا

علیرضا آخوندی

۹۷۳۱۱۰۷

سوال اول :

در تسک دسته بندی وظیفه مدل پیش بینی کردن برچسب برای یک داده ورودی است. در واقع خروجی این مدل ها گسسته است. در حالی که در رگرسیون خروجی پیوسته بوده و با توجه به ورودی مدل خروجی یک کمیت پیوسته و مقدار عددی است.

مثال های دسته بندی:

تشخیص بیماری: در این کاربرد، مدل تصاویر با فرمت های مختلف (MRI ، ...) دریافت کرده و مبتلا بوده یا عدم ابتلا به یک بیماری خاص را تشخیص می دهد.

تشخیص مانع در ماشین های خودران: در این کاربرد مدل یا با کمک داده های سنسور فاصله یاب یا با کمک تصویر گرفته شده از دوربین روبروی خودرو وجود یا عدم وجود یک مانع را تشخیص میدهد. به دست آوردن خصوصیات چهره : برای این که بتوانیم خصوصیات چهره یک فرد را (مانند رنگ پوست، لبخند، عصبانی یا خوشحال بودن و ...) را به دست آوریم برای هر ویژگی میتوانیم یک دسته بند ایجاد کنیم.

Authentication: برای این که با کمک تصویر صورت یک فرد تشخیص دهیم که آیا آن فرد اجازه دسترسی به منابعی را (برای مثال ورود به آزمایشگاه) دارد می توانیم از دسته بند ها استفاده کنیم. وفاداری: اینکه متوجه شویم مشتری یک شرکت خاص به آن شرکت وفادار است یا خیر بر اساس ویژگی های مختلف آن مشتری مانند: میزان موجودی، میزان خرید های آن مشتری، میزان استفاده آن مشتری از شرکت و...

مثال های رگرسیون:

پیش بینی سرعت خودرو: مدل با دریافت چند فریم از دوربین جلوی خودرو سرعت آن را پیش بینی میکند. کاربرد های پزشکی: با توجه به سوابق بیمار می توان مقدار قند خون یا هر ویژگی دیگری را پیشبینی کرد. معدل یک دانشجو: با توجه به سوابق دانشجو و روحیات او می توان معدل ترم او را پیش بینی کرد. میزان مصرف یک دارو: برای درمان یک بیماری، میتوان میزان مصرف یک دارو را برای هر بیمار با توجه به سوابق او تنظیم کرد.

سهام : پیش بینی نوسانات سهام با کمک داده های از قبل جمع شده.

سوال دوم: بخش اول:

Accuracy: دقت یک مدل برابر است با تعداد داده هایی که درست طبقه بندی شده اند تقسیم بر تعداد کل داده ها و فرمول آن نیز به صورت زیر است:

$$Accuracy = \frac{TN + TP}{TN + FP + TP + FN}$$

Recall: نشان می دهد چند درصد از داده هایی که واقعا مربوط به یک دسته هستند در آن دسته قرار گرفته اند که فرمول آن به صورت زیر است:

$$Recall = \frac{TP}{TP + FN}$$

Precision: همان دقت است که نسبت پاسخ های مثبت درست را به پاسخ های مثبت مشخص میکند که به صورت زیر محاسبه میشود:

$$Precision = \frac{TP}{TP + FP}$$

F1-Score: حاصل ضرب precision و recall تقسیم بر حاصل جمع آن ها ضرب در دو.

بخش دوم :

- سناریو اول: فرض کنید که می خواهیم عملکرد مدلی را ارزیابی کنیم که وظیفه آن تشخیص یک نوع سرطان است. فرض کنید که در مجموعه داده ما داده های مربوط به سوابق 1000 بیمار قرار دارد و تنها 10 مورد آنها این سرطان را دارند. در این صورت اگر از معیار accuracy استفاده کنیم آنگاه مدلی که به ازای تمام سوابق پاسخ منفی می دهد، دقت 99 درصد خواهد داشت که دقت خوبی است. اما در عمل این مدل قادر به شناسایی این سرطان نیست. برای ارزیابی عملکرد این مدل بهتر است که از precision و recall استفاده کنیم.
- سناریو دوم: حال فرض کنید می خواهیم عملکرد مدلی را ارزیابی کنیم که با توجه به سابقه کاربر پیش بینی می کند که آیا کاربر بر روی یک تبلیغ به خصوص کلیک می کند یا خیر. برخلاف سناریو پیشین، TN در این سناریو مقدار زیادی است و تاثیر آن را فقط در accuracy می توانیم ببینیم. بنابراین این معیار از دیگر معیار ها برای ارزیابی بهتر است.
- سناریو سوم: برای ارزیابی یک موتور جستجو معیار های accuracy و precision به اندازه recall برای این تسک مناسب نیستند. دلیل این ادعا این است که سایت های مرتبط با یک کوئری بسیار کمتر از سایت های غیر مرتبط است. و از انجایی که تمرکز recall تنها بر روی اسناد بازگردانده شده است، این معیار مناسب تر است.

سوال سوم:

مسئله سوم:

$$info(D) = - \sum_{i=1}^n p_i \log_2(p_i)$$

$$D_{\text{بیماری قلبی}} = [2, -2]$$

$$info(D) = - \frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$info_A(D_{\text{بیماری قلبی}}) = \sum_{j=1}^7 \frac{|D_j|}{|D|} \times info(D_j)$$

$$info_{\text{در سینه}}(D) = \frac{3}{4} \left(-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right) + \frac{1}{4} \times (1 \log_2 1) = 0.69$$

0.91

$$gain(\text{در سینه}) = 1 - 0.69 = 0.31$$

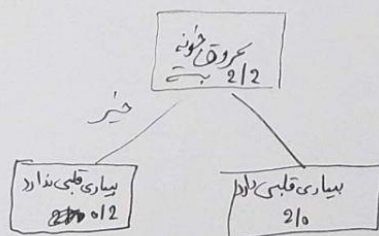
$$info_{\text{گردش خون}}(D) = \frac{2}{4} \times \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{2}{4} \times 1 = 1$$

$$gain(\text{گردش خون}) = 1 - 1 = 0$$

$$info_{\text{عروق خونی}}(D) = \frac{2}{4} \times 0 + \frac{2}{4} \times 1 = 0.5$$

$$gain(\text{عروق خونی}) = 1 - 0.5 = 0.5$$

باتوجه به این مای بدست آمده محقق می بیند بیشترین وزن را داشته و به عنوان ریشه انتخاب می شود.



سوال چهارم:

سوال چهارم:

$$info(D) = -\frac{3}{7} \log_2 \frac{3}{7} - \frac{4}{7} \log_2 \frac{4}{7} = 0.98$$

دوست داشتن سران
کلاه قرمزی

$$info(D) = \frac{4}{7} \left(-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) + \frac{3}{7} \left(-\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \right)$$

دوست داشتن
پاه کزن
= 0.85

$$gain(D) = 0.14$$

دوست داشتن
پاه کزن

$$info(D) = \frac{4}{7} \left(-\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} \right) + \frac{3}{7} \times 0 = 0.47$$

دوست داشتن
آب گازدار

$$gain(D) = 0.98 - 0.47 = 0.51$$

دوست داشتن
آب گازدار

از آن دسته که یک تغییر عددی است باید یک threshold برای آن انتخاب کنیم.

| | | | | | | | | |
|---|----|----|----|----|----|----|----|----|
| 7 | 12 | 15 | 18 | 35 | 38 | 44 | 50 | 83 |
| - | - | + | + | + | - | - | - | - |

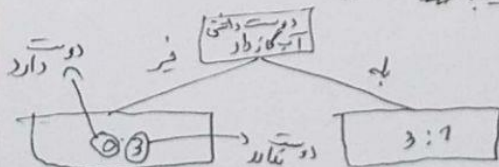
در دوسه 44 و 15 یک تغییر مشاهده داشتیم. برای مثال 15 را به عنوان مرز انتخاب می کنیم.

$$info(D) = \frac{2}{7} \times (-1 \log 1) + \frac{5}{7} \left(-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right) = 0.7$$

دوست داشتن
15 < 15

$$gain(D) = 0.29$$

بالاترین رتبه برای به دست آمده، درخت را به شکل زیر است:



ادامه 4:
برای شانه به به نیاز به گره دیگری داریم

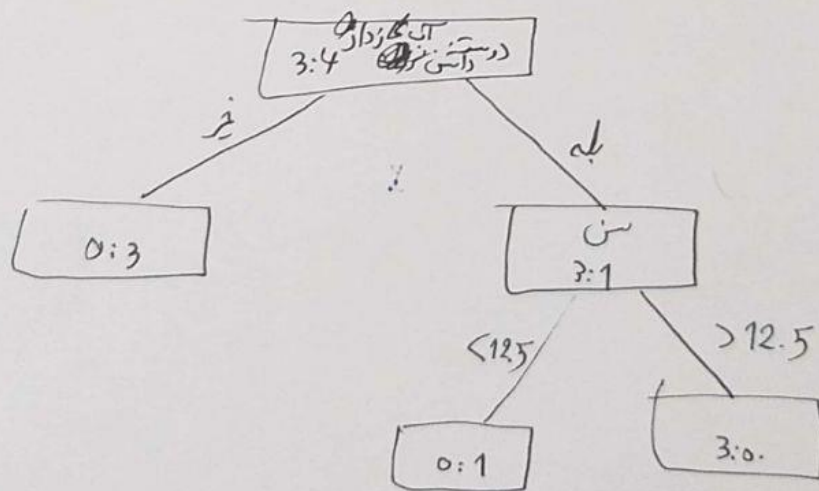
$$\text{info}_{\text{دانشگاه}} = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.81$$

$$\text{info}_{\text{پایه}} = \frac{2}{4} \left(-\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} \right) + \frac{2}{4} \times (-1 \log_2 1) = 0.5$$

$$\text{gain}(\text{دانشگاه} - \text{پایه}) = 0.81 - 0.5 = 0.31$$

$$\text{info}_{13} = \frac{1}{4} \times 0 + \frac{3}{4} \times 1 = 0.75$$

$$\text{gain}(\text{سن} > 13) = 0.81$$



سوال پنجم: شاخص جینی معیاری برای برای اندازه گیری ناهمگامی در یک است درجه ای که یک متغیر به اشتباه دسته بندی شده باشد را به دست می آورد. مقدار صفر نشان دهنده این است که تمامی عناصر به یک کلاسی تعلق دارند در حالی که مقدار 1 نشان دهنده این است که تمامی عناصر به طور رندوم بین کلاسی های مختلف پخش شده اند.

$$Gini = 1 - \sum_{i=1}^n (p_i)^2$$

بر نشان دهنده احتمال این که یک شیء خاص در کلاس خاصی قرار گیرد. برای نشان درخت سوال چهارم را در نظر بگیرید.

$$P(\text{آب گزدار} = +) = \frac{4}{7} \quad P(\text{آب گزدار} = -) = \frac{3}{7}$$

$$P(\text{دریا کلاه قرمز} = + \text{ and } \text{آب گزدار} = +) = \frac{3}{7}$$

$$P(\text{دریا کلاه قرمز} = + \mid \text{آب گزدار} = +) = \frac{3}{4}$$

$$P(\text{دریا کلاه قرمز} = - \mid \text{آب گزدار} = +) = \frac{1}{4}$$

$$Gini_{\text{آب گزدار}} = 1 - \left(\frac{4}{7} + \frac{1}{7} \right) = \frac{3}{7}$$

$$P(\text{دریا کلاه قرمز} = + \mid \text{آب گزدار} = -) = 0$$

$$P(\text{دریا کلاه قرمز} = - \mid \text{آب گزدار} = -) = 1$$

$$Gini_{\text{آب گزدار} = -} = 1 - (0^2 + 1^2) = 0$$

$$Gini(\text{آب گزدار}) = P(\text{آب گزدار} = +) \times Gini(\text{آب گزدار} = +) +$$

$$P(\text{آب گزدار} = -) \times Gini(\text{آب گزدار} = -) = \frac{4}{7} \times \frac{3}{7} + \frac{3}{7} \times 0 = \frac{12}{49}$$

سوال ششم:

Overfitting زمانی رخ میدهد که مدل ما بر روی داده های آموزشی بیش از اندازه آموزش داده شود و یا مدلی که برای تسک طراحی شده است از پیچیدگی بالایی برخوردار است. نتیجه آن دقت بالا در داده آموزش و دقت پایین در داده تست است. برای مثال انتخاب یک شبکه بسیار عمیق برای یک تسک ساده می تواند به این مشکل بخورد یا قرار ندادن **limit** بر روی عمق یک درخت تصمیم گیری نیز می تواند این مشکل را ایجاد کند.

برای جلوگیری از این مشکل می توان کار های زیر را انجام داد:

افزایش داده: از دلایل این مشکل داده کم است. افزایش تنوع در مجموعه داده آموزش می تواند این مشکل را تا حد خوبی برطرف کند.

ساده سازی مدل: راه دیگر ساده سازی ساختار مدل مورد استفاده است. برای مثال اگر این مدل یک شبکه عصبی است، می توان تعداد لایه های آن را کاهش داد یا اگر یک درخت تصمیم گیری است می توان عمق آن را کاهش داد.

Dropout: در این روش با در نظر گرفتن یک احتمال برای هر لایه از یک شبکه عصبی، برخی از لایه ها را در زمان آموزش نادیده می گیریم.

Regularization: تجربه ثابت کرده است که هر چه قدر اندازه وزن های یک مدل کوچکتر باشد، آن مدل ساده تر است. برای همین می توان از مکانیزمی استفاده کرد که در هنگام آموزش، مدل را برای داشتن وزن های بزرگ تنبیه کند. از این روش ها می توان به l_1 و l_2 regularization اشاره کرد.

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^N |w_i|$$

Loss function with L1 regularisation

$$Loss = Error(y, \hat{y}) + \lambda \sum_{i=1}^N w_i^2$$

Loss function with L2 regularisation